

One-Class Support Measure Machines for Group Anomaly Detection

Krikamol Muandet

Empirical Inference Department
Max Planck Institute for Intelligent Systems
Spemannstraße 38, 72076 Tübingen

Bernhard Schölkopf

Empirical Inference Department
Max Planck Institute for Intelligent Systems
Spemannstraße 38, 72076 Tübingen

Abstract

We propose one-class support measure machines (OCSMMs) for group anomaly detection. Unlike traditional anomaly detection, OCSMMs aim at recognizing anomalous aggregate behaviors of data points. The OCSMMs generalize well-known one-class support vector machines (OCSVMs) to a space of probability measures. By formulating the problem as quantile estimation on distributions, we can establish interesting connections to the OCSVMs and variable kernel density estimators (VKDEs) over the input space on which the distributions are defined, bridging the gap between large-margin methods and kernel density estimators. In particular, we show that various types of VKDEs can be considered as solutions to a class of regularization problems studied in this paper. Experiments on Sloan Digital Sky Survey dataset and High Energy Particle Physics dataset demonstrate the benefits of the proposed framework in real-world applications.

1 Introduction

Anomaly detection is one of the most important tools in all data-driven scientific disciplines. Data that do not conform to the expected behaviors often bear some interesting characteristics and can help domain experts better understand the problem at hand. However, in the era of data explosion, the anomaly may appear not only in the data themselves, but also as a result of their interactions. The main objective of this paper is to investigate the latter type of anomalies. To be consistent with the previous works (Póczos et al. 2011, Xiong et al. 2011a;b), we will refer to this problem as a group anomaly detection, as opposed to a traditional point anomaly detection.

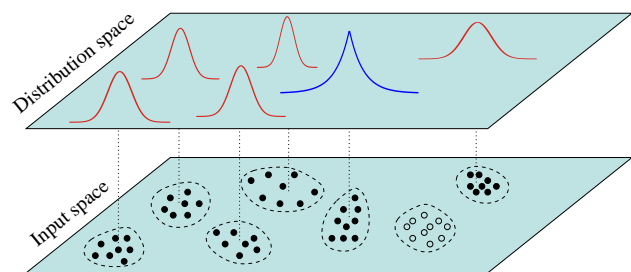


Figure 1: An illustration of two types of group anomalies. An anomalous group may be a group of anomalous samples which is easy to detect (unfilled points). In this paper, we are interested in detecting anomalous groups of normal samples (filled points) which is more difficult to detect because of the higher-order statistics. Note that group anomaly we are interested in can only be observed in the space of distributions.

Like traditional point anomaly detection, the group anomaly detection refers to a problem of finding patterns in groups of data that do not conform to expected behaviors (Póczos et al. 2011, Xiong et al. 2011a;b). That is, an ultimate goal is to detect interesting aggregate behaviors of data points among several groups. In principle, anomalous groups may consist of individually anomalous points, which are relatively easy to detect. On the other hand, anomalous groups of relatively normal points, whose behavior as a group is unusual, is much more difficult to detect. In this work, we are interested in the latter type of group anomalies. Figure 1 illustrates this scenario.

Group anomaly detection may shed light in a wide range of applications. For example, a Sloan Digital Sky Survey (SDSS) has produced a tremendous amount of astronomical data. It is therefore very crucial to detect rare objects such as stars, galaxies, or quasars that might lead to a scientific discovery. In addition to individual celestial objects, investigating groups of them may help astronomers understand the universe on larger scales. For instance, the anomalous

group of galaxies, which is the smallest aggregates of galaxies, may reveal interesting phenomena, e.g., the gravitational interactions of galaxies.

Likewise, a new physical phenomena in high energy particle physics such as Higgs boson appear as a tiny excesses of certain types of collision events among a vast background of known physics in particle detectors (Bhat 2011, Vatanen et al. 2012). Investigating each collision event individually is no longer sufficient as the individual events may not be anomalies by themselves, but their occurrence together as a group is anomalous. Hence, we need a powerful algorithm to detect such a rare and highly structured anomaly.

Lastly, the algorithm proposed in this paper can be applied to point anomaly detection with substantial and heterogeneous uncertainties. For example, it is often costly and time-consuming to obtain the full spectra of astronomical objects. Instead, relatively noisier measurements are usually made. In addition, the estimated uncertainty which represents the uncertainty one would obtain from multiple observations is also available. Incorporating these uncertainties has been shown to improve the performance of the learning systems (Bovy et al. 2011, Kirkpatrick et al. 2011, Ross et al. 2012).

The anomaly detection has been intensively studied (Chandola et al. (2009) and references therein). However, few attempts have been made on developing successful group anomaly detection algorithms. For example, a straightforward approach is to define a set of features for each group and apply standard point anomaly detection (Chan and Mahoney 2005). Despite its simplicity, this approach requires a specific domain knowledge to construct appropriate sets of features. Another possibility is to first identify the individually anomalous points and then find their aggregations (Das et al. 2008). Again, this approach relies only on the detection of anomalous points and thus cannot find the anomalous groups in which their members are perfectly normal. Successful group anomaly detectors should be able to incorporate the higher-order statistics of the groups.

Recently, a family of hierarchical probabilistic models based on a Latent Dirichlet Allocation (LDA) (Blei et al. 2003) has been proposed to cope with both types of group anomalies (Xiong et al. 2011a;b). In these models, the data points in each group are assumed to be one of the K different types and generated by a mixture of K Gaussian distributions. Although the distributions over these K types can vary across M groups, they share common generator. The groups that have small probabilities under the model are marked as anomalies using scoring criteria defined

as a combination of a point-based anomaly score and a group-based anomaly score. The Flexible Genre Model (FGM) recently extends this idea to model more complex group structures (Xiong et al. 2011a).

Instead of employing a generative approach, we propose a simple and efficient discriminative way of detecting group anomaly. In this work, M groups of data points are represented by a set of M probability distributions assumed to be i.i.d. realization of some unknown distribution \mathcal{P} . In practice, only i.i.d samples from these distributions are observed. Hence, we can treat group anomaly detection as detecting the anomalous distributions based on their empirical samples. To allow for a practical algorithm, the distributions are mapped into the reproducing kernel Hilbert space (RKHS) using the kernel mean embedding. By working directly with the distributions, the higher-order information arising from the aggregate behaviors of the data points can be incorporated efficiently.

2 Quantile Estimation on Probability Distributions

Let \mathcal{X} denote a non-empty input space with associated σ -algebra \mathcal{A} , \mathbb{P} denote the probability distribution on $(\mathcal{X}, \mathcal{A})$, and $\mathfrak{P}_{\mathcal{X}}$ denote the set of all probability distributions on $(\mathcal{X}, \mathcal{A})$. The space $\mathfrak{P}_{\mathcal{X}}$ is endowed with the topology of weak convergence and the associated Borel σ -algebra.

We assume that there exists a distribution \mathcal{P} on $\mathfrak{P}_{\mathcal{X}}$, where $\mathbb{P}_1, \dots, \mathbb{P}_{\ell}$ are i.i.d. realizations from \mathcal{P} , and the sample S_i is made of n_i i.i.d. samples distributed according to the distribution \mathbb{P}_i . In this work, we observe ℓ samples $S_i = \{x_k^{(i)}\}_{1 \leq k \leq n_i}$ for $i = 1, \dots, \ell$. For each sample S_i , $\hat{\mathbb{P}}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} \delta_{x_j^{(i)}}$ is the associated empirical distribution of \mathbb{P}_i .

In this work, we formulate a group anomaly detection problem as learning quantile function $q : \mathfrak{P}_{\mathcal{X}} \rightarrow \mathbb{R}$ to estimate the support of \mathcal{P} . Let \mathcal{C} be a class of measurable subsets of $\mathfrak{P}_{\mathcal{X}}$ and λ be a real-valued function defined on \mathcal{C} , the quantile function w.r.t. $(\mathcal{P}, \mathcal{C}, \lambda)$ is

$$q(\beta) = \inf\{\lambda(C) : \mathcal{P}(C) \geq \beta, C \in \mathcal{C}\},$$

where $0 < \beta \leq 1$. In this paper, we consider when λ is Lebesgue measure, in which case $C(\beta)$ is the minimum volume $C \in \mathcal{C}$ that contains at least a fraction β of the probability mass of \mathcal{P} . Thus, the function q can be used to test if any test distribution \mathbb{P}_t is anomalous w.r.t. the training distributions.

Rather than estimating $C(\beta)$ in the space of distributions directly, we first map the distributions into a feature space via a positive semi-definite kernel k .

Our class \mathcal{C} is then implicitly defined as the set of half-spaces in the feature space. Specifically, $C_{\mathbf{w}} = \{\mathbb{P} \mid f_{\mathbf{w}}(\mathbb{P}) \geq \rho\}$ where (\mathbf{w}, ρ) are respectively a weight vector and an offset parametrizing a hyperplane in the feature space associated with the kernel k . The optimal (\mathbf{w}, ρ) is obtained by minimizing a regularizer which controls the smoothness of the estimated function describing C .

3 One-Class Support Measure Machines

In order to work with the probability distributions efficiently, we represent the distributions as mean functions in a reproducing kernel Hilbert space (RKHS) (Berlinet and Agnan 2004, Smola et al. 2007). Formally, let \mathcal{H} denote an RKHS of functions $f : \mathcal{X} \rightarrow \mathbb{R}$ with reproducing kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. The kernel mean map from $\mathfrak{P}_{\mathcal{X}}$ into \mathcal{H} is defined as

$$\mu : \mathfrak{P}_{\mathcal{X}} \rightarrow \mathcal{H}, \quad \mathbb{P} \mapsto \int_{\mathcal{X}} k(x, \cdot) d\mathbb{P}(x) . \quad (1)$$

We assume that $k(x, \cdot)$ is bounded for any $x \in \mathcal{X}$. For any \mathbb{P} , letting $\mu_{\mathbb{P}} = \mu(\mathbb{P})$, one can show that $\mathbb{E}_{\mathbb{P}}[f] = \langle \mu_{\mathbb{P}}, f \rangle_{\mathcal{H}}$, for all $f \in \mathcal{H}$.

The following theorem due to Fukumizu et al. (2004) and Sriperumbudur et al. (2010) gives a promising property of representing distributions as mean elements in the RKHS.

Theorem 1. *The kernel k is characteristic if and only if the map (1) is injective.*

Examples of characteristic kernels include Gaussian RBF kernel and Laplace kernel. Using the characteristic kernel k , Theorem 1 implies that the map (1) preserves all information about the distributions. Hence, one can apply many existing kernel-based learning algorithms to the distributions as if they are individual samples with no information loss.

Intuitively, one may view the mean embeddings of the distributions as their feature representations. Thus, our approach is in line with previous attempts in group anomaly detection that find a set of appropriate features for each group. On the one hand, however, the mean embedding approach captures all necessary information about the groups without relying heavily on a specific domain knowledge. On the other hand, it is flexible to choose the feature representation that is suitable to the problem at hand via the choice of the kernel k .

3.1 OCSMM Formulation

Using the mean embedding representation (1), the primal optimization problem for one-class SMM can be

subsequently formulated in an analogous way to the one-class SVM (Schölkopf et al. 2001) as follow:

$$\underset{\mathbf{w}, b, \xi, \rho}{\text{minimize}} \quad \frac{1}{2} \langle \mathbf{w}, \mathbf{w} \rangle_{\mathcal{H}} - \rho + \frac{1}{\nu \ell} \sum_{i=1}^{\ell} \xi_i \quad (2a)$$

$$\text{subject to} \quad \langle \mathbf{w}, \mu_{\mathbb{P}_i} \rangle_{\mathcal{H}} \geq \rho - \xi_i, \xi_i \geq 0 \quad (2b)$$

where ξ_i denote slack variables and $\nu \in (0, 1]$ is a trade-off parameter corresponding to an expected fraction of outliers within the feature space. The trade-off ν is an upper bound on the fraction of outliers and lower bound on the fraction of support measures (Schölkopf et al. 2001).

The trade-off parameter ν plays an important role in group anomaly detection. Small ν implies that anomalous groups are rare compared to the normal groups. Too small ν leads to some anomalous groups being rejected. On the other hand, large ν implies that anomalous groups are common. Too large ν leads to some normal groups being accepted as anomaly. As group anomaly is subtle, one need to choose ν very carefully to reduce the effort in the interpretation of the results.

By introducing Lagrange multipliers α , we have $\mathbf{w} = \sum_{i=1}^{\ell} \alpha_i \mu_{\mathbb{P}_i} = \sum_{i=1}^{\ell} \alpha_i \mathbb{E}_{\mathbb{P}_i}[k(x, \cdot)]$ and the dual form of (2) can be written as

$$\underset{\alpha}{\text{minimize}} \quad \frac{1}{2} \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} \alpha_i \alpha_j \langle \mu_{\mathbb{P}_i}, \mu_{\mathbb{P}_j} \rangle_{\mathcal{H}} \quad (3a)$$

$$\text{subject to} \quad 0 \leq \alpha_i \leq \frac{1}{\nu \ell}, \sum_{i=1}^{\ell} \alpha_i = 1 . \quad (3b)$$

Note that the dual form is a quadratic programming and depends on the inner product $\langle \mu_{\mathbb{P}_i}, \mu_{\mathbb{P}_j} \rangle_{\mathcal{H}}$. Given that we can compute $\langle \mu_{\mathbb{P}_i}, \mu_{\mathbb{P}_j} \rangle_{\mathcal{H}}$, we can employ the standard QP solvers to solve (3).

3.2 Kernels on Probability Distributions

From (3), we can see that $\mu_{\mathbb{P}}$ is a feature map associated with the kernel $K : \mathfrak{P}_{\mathcal{X}} \times \mathfrak{P}_{\mathcal{X}} \rightarrow \mathbb{R}$, defined as $K(\mathbb{P}_i, \mathbb{P}_j) = \langle \mu_{\mathbb{P}_i}, \mu_{\mathbb{P}_j} \rangle_{\mathcal{H}}$. It follows from Fubini's theorem and reproducing property of \mathcal{H} that

$$\begin{aligned} \langle \mu_{\mathbb{P}_i}, \mu_{\mathbb{P}_j} \rangle_{\mathcal{H}} &= \iint \langle k(x, \cdot), k(y, \cdot) \rangle_{\mathcal{H}} d\mathbb{P}_i(x) d\mathbb{P}_j(y) \\ &= \iint k(x, y) d\mathbb{P}_i(x) d\mathbb{P}_j(y) . \end{aligned} \quad (4)$$

Hence, K is a positive definite kernel on $\mathfrak{P}_{\mathcal{X}}$. Given the sample sets S_1, \dots, S_{ℓ} , one can estimate (4) by

$$K(\widehat{\mathbb{P}}_i, \widehat{\mathbb{P}}_j) = \frac{1}{n_i \cdot n_j} \sum_{k=1}^{n_i} \sum_{l=1}^{n_j} k(x_k^{(i)}, x_l^{(j)}) \quad (5)$$

where $x_k^{(i)} \in S_i$, $x_l^{(j)} \in S_j$, and n_i is the number of samples in S_i for $i = 1, \dots, \ell$.

Previous works in kernel-based anomaly detection have shown that the Gaussian RBF kernel is more suitable than some other kernels such as polynomial kernels (Hoffmann 2007). Thus we will focus primarily on the Gaussian RBF kernel given by

$$k_\sigma(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right), \quad x, x' \in \mathcal{X} \quad (6)$$

where $\sigma > 0$ is a bandwidth parameter. In the sequel, we denote the reproducing kernel Hilbert space associated with kernel k_σ by \mathcal{H}_σ . Also, let $\Phi : \mathcal{X} \rightarrow \mathcal{H}_\sigma$ be a feature map such that $k_\sigma(x, x') = \langle \Phi(x), \Phi(x') \rangle_{\mathcal{H}_\sigma}$.

In group anomaly detection, we always observe the i.i.d. samples from the distribution underlying the group. Thus, it is natural to use the empirical kernel (5). However, one may relax this assumption and apply the kernel (4) directly. For instance, if we have a Gaussian distribution $\mathbb{P}_i = \mathcal{N}(m_i, \Sigma_i)$ and a Gaussian RBF kernel k_σ , we can compute the kernel analytically by

$$K(\mathbb{P}_i, \mathbb{P}_j) = \frac{\exp\left(-\frac{1}{2}(m_i - m_j)^\top B^{-1}(m_i - m_j)\right)}{\left|\frac{1}{\sigma^2}\Sigma_i + \frac{1}{\sigma^2}\Sigma_j + \mathbf{I}\right|^{\frac{1}{2}}} \quad (7)$$

where $B = \Sigma_i + \Sigma_j + \sigma^2\mathbf{I}$. This kernel is particularly useful when one wants to incorporate the point-wise uncertainty of the observation into the learning algorithm (Muandet et al. 2012). More details will be given in Section 4.2 and 5).

4 Theoretical Analysis

This section presents some theoretical analyses. The geometrical interpretation of OCSMMs is given in Section 4.1. Then, we discuss the connection of OCSMM to the kernel density estimator in Section 4.2. In the sequel, we will focus on the translation-invariant kernel function to simplify the analysis.

4.1 Geometric Interpretation

For translation-invariant kernel, $k(x, x)$ is constant for all $x \in \mathcal{X}$. That is, $\|\Phi(x)\|_{\mathcal{H}} = \tau$ for some constant ρ . This implies that all of the images $\Phi(x)$ lie on the sphere in the feature space (cf. Figure 2a). Consequently, the following inequality holds

$$\|\mu_{\mathbb{P}}\|_{\mathcal{H}} = \left\| \int k(x, \cdot) d\mathbb{P}(x) \right\|_{\mathcal{H}} \leq \int \|k(x, \cdot)\|_{\mathcal{H}} d\mathbb{P}(x) = \tau,$$

which shows that all mean embeddings lie inside the sphere (cf. Figure 2a). As a result, we can establish the existence and uniqueness of the separating hyperplane \mathbf{w} in (2) through the following theorem.

Theorem 2. *There exists a unique separating hyperplane \mathbf{w} as a solution to (2) that separates $\mu_{\mathbb{P}_1}, \mu_{\mathbb{P}_2}, \dots, \mu_{\mathbb{P}_\ell}$ from the origin.*

Proof. Due to the separability of the feature maps $\Phi(x)$, the convex hull of the mean embeddings $\mu_{\mathbb{P}_1}, \mu_{\mathbb{P}_2}, \dots, \mu_{\mathbb{P}_\ell}$ does not contain the origin. The existence and uniqueness of the hyperplane then follows from the supporting hyperplane theorem (Schölkopf and Smola 2001). ■

By Theorem 2, the OCSMM is a simple generalization of OCSVM to the space of probability distributions. Furthermore, the straightforward generalization will allow for a direct application of an efficient learning algorithm as well as existing theoretical results.

There is a well-known connection between the solution of OCSVM with translation invariant kernels and the center of the minimum enclosing sphere (MES) (Tax and Duin 1999; 2004). Intuitively, this is not the case for OCSMM, even when the kernel k is translation-invariant, as illustrated in Figure 2b. Fortunately, the connection between OCSMM and MES can be made precise by applying the spherical normalization

$$\langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}} \mapsto \frac{\langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}}}{\sqrt{\langle \mu_{\mathbb{P}}, \mu_{\mathbb{P}} \rangle_{\mathcal{H}} \langle \mu_{\mathbb{Q}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}}}} \quad (8)$$

After the normalization, $\|\mu_{\mathbb{P}}\|_{\mathcal{H}} = 1$ for all $\mathbb{P} \in \mathfrak{P}_{\mathcal{X}}$. That is, all mean embeddings lie on the unit sphere in the feature space. Consequently, the OCSMM and MES are equivalent after the normalization.

Given the equivalence between OCSMM and MES, it is natural to ask if the spherical normalization (8) preserves the injectivity of the Hilbert space embedding. In other words, is there an information loss after the normalization? The following theorem answers this question for kernel k that satisfies some reasonable assumptions.

Theorem 3. *Assume that k is characteristic and the samples are linearly independent in the feature space \mathcal{H} . Then, the spherical normalization preserves the injectivity of the mapping $\mu : \mathfrak{P}_{\mathcal{X}} \rightarrow \mathcal{H}$.*

Proof. Let us assume the normalization does not preserve the injectivity of the mapping. Thus, there exist two distinct probability distributions \mathbb{P} and \mathbb{Q} for which

$$\begin{aligned} \mu_{\mathbb{P}} &= \mu_{\mathbb{Q}} \\ \int k(x, \cdot) d\mathbb{P}(x) &= \int k(x, \cdot) d\mathbb{Q}(x) \\ \int k(x, \cdot) d(\mathbb{P} - \mathbb{Q})(x) &= 0. \end{aligned}$$

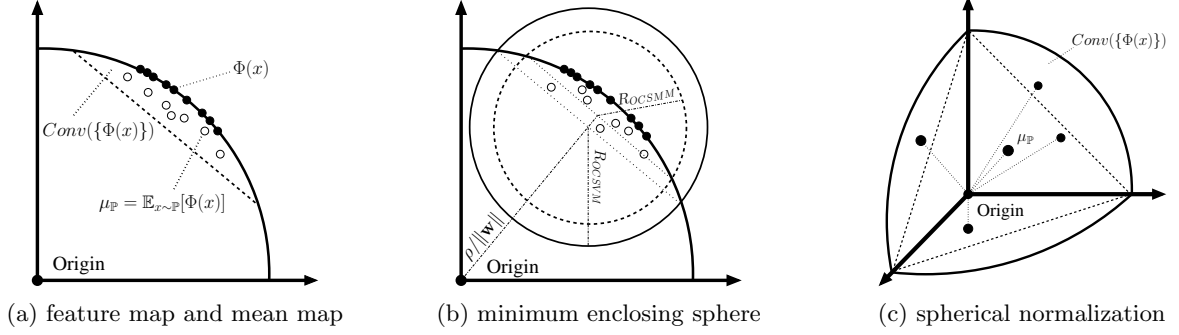


Figure 2: (a) The two dimensional representation of the RKHS of Gaussian RBF kernels. Since the kernels depend only on $x - x'$, $k(x, x')$ is constant. Therefore, all feature maps $\Phi(x)$ (black dots) lie on a sphere in feature space. Hence, for any probability distribution \mathbb{P} , its mean embedding $\mu_{\mathbb{P}}$ always lies in the convex hull of the feature maps, which in this case, forms a segment of the sphere. (b) In general, the solution of OCSMM is different from the minimum enclosing sphere. (c) Three dimensional sphere in the feature space. For the Gaussian RBF kernel, the kernel mean embeddings of all distributions always lie inside the segment of the sphere. In addition, the angle between any pair of mean embeddings is always greater than zero. Consequently, the mean embeddings can be scaled, e.g., to lie on the sphere, and the map is still injective.

As $\mathbb{P} \neq \mathbb{Q}$, the last equality holds if and only if there exists $x \in \mathcal{X}$ for which $k(x, \cdot)$ are linearly dependent, which contradicts the assumption. Consequently, the spherical normalization must preserve the injectivity of the mapping. ■

The Gaussian RBF kernel satisfies the assumption given in Theorem 3 as the kernel matrix will be full-rank and thereby the samples are linearly independent in the feature space. Figure 2c depicts an effect of the spherical normalization.

It is important to note that the spherical normalization does not necessarily improve the performance of the OCSMM. It ensures that all the information about the distributions are preserved.

4.2 OCSMM and Density Estimation

In this section we make a connection between the OCSMM and kernel density estimation (KDE). First, we give a definition of the KDE. Let x_1, x_2, \dots, x_n be an i.i.d. samples from some distribution F with unknown density f , the KDE of f is defined as

$$\hat{f}(y) = \frac{1}{nh} \sum_{i=1}^n k\left(\frac{y - x_i}{h}\right) \quad (9)$$

For \hat{f} to be a density, we require that the kernel satisfies $k(\cdot, \cdot) \geq 0$ and $\int k(x, \cdot) dx = 1$, which includes, for example, the Gaussian kernel, the multivariate Student kernel, and the Laplacian kernel.

When $\nu = 1$, it is well-known that, under some technical assumptions, the OCSVM corresponds exactly to the KDE (Schölkopf et al. 2001). That is, the solution

\mathbf{w} of (2) can be written as a uniform sum over training samples similar to (9). Moreover, setting $\nu < 1$ yields a sparse representation where the summand consists of only support vectors of the OCSVM.

Interestingly, we can make a similar correspondence between the KDE and the OCSMM. It follows from Lemma 4 of Muandet et al. (2012) that for certain classes of training probability distributions, the OCSMM on these distributions corresponds to the OCSVM on some training samples equipped with an appropriate kernel function. To understand this connection, consider the OCSMM with the Gaussian RBF kernel k_{σ} and isotropic Gaussian distributions $\mathcal{N}(m_1; \sigma_1^2), \mathcal{N}(m_2; \sigma_2^2), \dots, \mathcal{N}(m_n; \sigma_n^2)$ ¹. We analyze this scenario under two conditions:

(C1) Identical bandwidth. If $\sigma_i = \sigma_j$ for all $1 \leq i, j \leq n$, the OCSMM is equivalent to the OCSVM on the training samples m_1, m_2, \dots, m_n with Gaussian RBF kernel $k_{\sigma^2 + \sigma_i^2}$ (cf. the kernel (7)). Hence, the OCSMM corresponds to the OCSVM on the means of the distributions with kernel of larger bandwidth.

(C2) Variable bandwidth. Similarly, if $\sigma_i \neq \sigma_j$ for some $1 \leq i, j \leq n$, the OCSMM is equivalent to the OCSVM on the training samples m_1, m_2, \dots, m_n with Gaussian RBF kernel $k_{\sigma^2 + \sigma_i^2}$. Note that the kernel bandwidth may be different at each training samples. Thus, OCSMM in this case corresponds to the OCSVM with variable bandwidth parameters.

¹We adopt the Gaussian distributions here for the sake of simplicity. More general statement for non-Gaussian distributions follows straightforwardly.

On the one hand, the above scenario allows the OCSVM to cope with noisy/uncertain inputs, leading to more robust point anomaly detection algorithm. That is, we can treat the means as the measurements and the covariances as the measurement uncertainties (cf. Section 5.2). On the other hand, one can also interpret the OCSMM when $\nu = 1$ as a generalization of traditional KDE, where we have a data-dependent bandwidth at each data point. This type of KDE is known in the statistics as variable kernel density estimators (VKDEs) (Abramson 1982, Breiman et al. 1977, Terrell and Scott 1992). For $\nu < 1$, the OCSMM gives a sparse representation of the VKDE.

Formally, the VKDE is characterized by (9) with an adaptive bandwidth $h(x_i)$. For example, the bandwidth is adapted to be larger where the data are less dense, with the aim to reduce the bias. There are basically two different views of VKDE. The first is known as a *balloon estimator* (Terrell and Scott 1992). Essentially, its bandwidth may depend only on the point at which the estimate is taken, i.e., the bandwidth in (9) may be written as $h(y)$. The second type of VKDE is a *sample smoothing estimator* (Terrell and Scott 1992). As opposed to the balloon estimator, it is a mixture of individually scaled kernels centered at each observation, i.e., the bandwidth is $h(x_i)$. The advantage of balloon estimator is that it has a straightforward asymptotic analysis, but the final estimator may not be a density. The sample smoothing estimator is a density if k is a density, but exhibits *non-locality*.

Both types of the VKDEs may be seen from the OCSMM point of view. Firstly, under the condition **(C1)**, the balloon estimator can be recovered by considering different test distribution $\mathbb{P}_t = \mathcal{N}(m_t; \sigma_t)$. As $\sigma_t \rightarrow 0$, one obtain the standard KDE on m_t . Similarly, the OCSMM under the condition **(C2)** with $\mathbb{P}_t = \delta_{m_t}$ gives the sample smoothing estimator. Interestingly, the OCSMM under the condition **(C2)** with $\mathbb{P}_t = \mathcal{N}(m_t; \sigma_t)$ results in a combination of these two types of the VKDEs.

In summary, we show that many variants of KDE can be seen as solutions to the regularization functional (2), and thereby provides an insight into a connection between large-margin approach and kernel density estimation.

5 Experiments

We firstly illustrate a fundamental difference between point and group anomaly detection problems. Then, we demonstrate an advantage of OCSMM on uncertain data when the noise is observed explicitly. Lastly, we compare the OCSMM with existing group anomaly detection techniques, namely, K -nearest neighbor

(KNN) based anomaly detection (Zhao and Saligrama 2009) with NP- L_2 divergence and NP-Renyi divergence (Póczos et al. 2011), and Multinomial Genre Model (MGM) (Xiong et al. 2011b) on Sloan Digital Sky Survey (SDSS) dataset and High Energy Particle Physics dataset.

Model Selection and Setup. One of the long-standing problems of one-class algorithms is model selection. Since no labeled data is available during training, we cannot perform cross validation. To encourage a fair comparison of different algorithms in our experiments, we will try out different parameter settings and report the best performance of each algorithm. We believe this simple approach should serve its purpose at reflecting the relative performance of different algorithms. We will employ the Gaussian RBF kernel (6) throughout the experiments. For the OCSVM and the OCSMM, the bandwidth parameter σ^2 is fixed at $\text{median}\{\|x_k^{(i)} - x_l^{(j)}\|^2\}$ for all i, j, k, l where $x_k^{(i)}$ denotes the k -th data point in the i -th group, and we consider $\nu = (0.1, 0.2, \dots, 0.9)$. The OCSVM treats group means as training samples. For synthetic experiments with OCSMM, we use the empirical kernel (5), whereas the non-linear kernel $K(\mathbb{P}_i, \mathbb{P}_j) = \exp(\|\mu_{\mathbb{P}_i} - \mu_{\mathbb{P}_j}\|_{\mathcal{H}}^2 / 2\gamma^2)$ will be used for real data where we set $\gamma = \sigma$. Our experiments suggest that these choices of parameters usually work well in practice. For KNN- L_2 and KNN-Renyi ($\alpha=0.99$), we consider when there are 3, 5, 7, 9, and 11 nearest neighbors. For MGM, we follow the same experimental setup as in Xiong et al. (2011b).

5.1 Synthetic Data

To illustrate the difference between point anomaly and group anomaly, we represent the group of data points by the 2-dimensional Gaussian distribution. We generate 20 normal groups with the covariance $\Sigma = [0.01, 0.008; 0.008, 0.01]$. The means of these groups are drawn uniformly from $[0, 1]$. Then, we generate 2 anomalous groups of Gaussian distributions whose covariances are rotated by 60 degree from the covariance Σ . Furthermore, we perturb one of the normal groups to make it relatively far from the rest of the dataset to introduce an additional degree of anomaly (cf. Figure 3a). Lastly, we generate 100 samples from each of these distributions to form the training set.

For the OCSVM, we represent each group by its empirical average. Since the expected proportion of outliers in the dataset is approximately 10%, we use $\nu = 0.1$ accordingly for both OCSVM and OCSMM. Figure 3a depicts the result which demonstrates that the OCSMM can detect anomalous aggregate patterns undetected by the OCSVM.

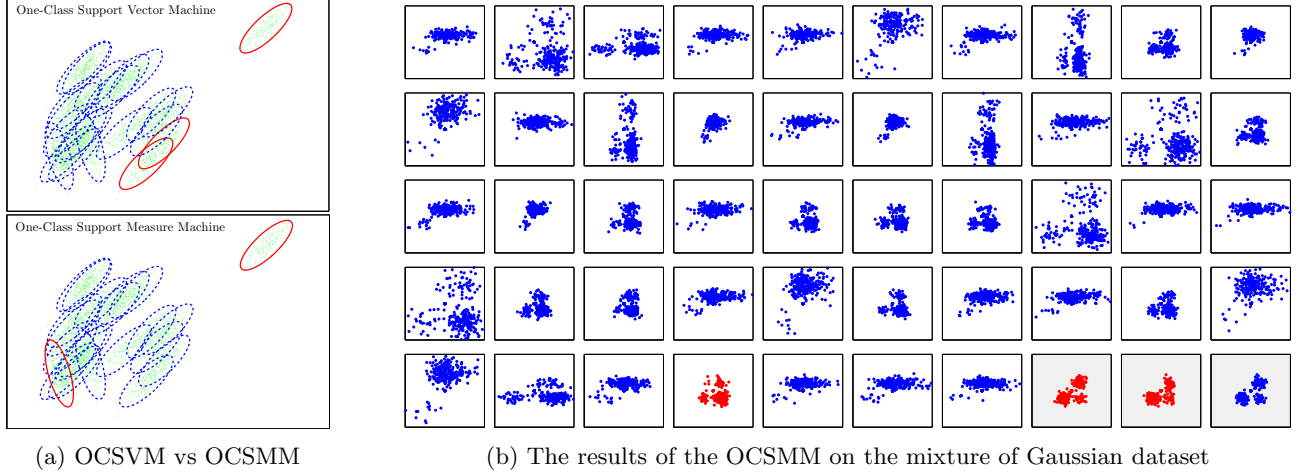


Figure 3: (a) The results of group anomaly detection on synthetic data obtained from the OCSVM and the OCSMM. Blue dashed ovals represent the normal groups, whereas red ovals represent the detected anomalous groups. The OCSVM is only able to detect the anomalous groups that are spatially far from the rest in the dataset, whereas the OCSMM also takes into account other higher-order statistics and therefore can also detect anomalous groups which possess distinctive properties. (b) The results of the OCSMM on the synthetic data of the mixture of Gaussian. The shaded boxes represent the anomalous groups that have different mixing proportion to the rest of the dataset. The OCSMM is able to detect the anomalous groups although they look reasonably normal and cannot be easily distinguished from other groups in the data set based only on an inspection.

Then, we conduct similar experiment as that in Xiong et al. (2011b). That is, the groups are represented as a mixture of four 2-dimensional Gaussian distributions. The means of the mixture components are $[-1, -1]$, $[1, -1]$, $[0, 1]$, $[1, 1]$ and the covariances are all $\Sigma = 0.15 \times \mathbf{I}_2$, where \mathbf{I}_2 denotes the 2D identity matrix. Then, we design two types of normal groups, which are specified by two mixing proportions $[0.22, 0.64, 0.03, 0.11]$ and $[0.22, 0.03, 0.64, 0.11]$, respectively. To generate a normal group, we first decide with probability $[0.48, 0.52]$ which mixing proportion will be used. Then, the data points are generated from mixture of Gaussian using the specified mixing proportion. The mixing proportion of the anomalous group is $[0.61, 0.1, 0.06, 0.23]$.

We generated 47 normal groups with $n_i \sim \text{Poisson}(300)$ instances in each group. Note that the individual samples in each group are perfectly normal compared to other samples. To test the performance of our technique, we inject the group anomalies, where the individual points are normal, but they together as a group look anomalous. In this anomalous group the individual points are samples from one of the $K = 4$ normal topics, but the mixing proportion was different from both of the normal mixing proportions. We inject 3 anomalous groups into the data set. The OCSMM is trained using the same setting as in the previous experiment. The results are depicted in Figure 3b.

5.2 Noisy Data

As discussed at the end of Section 3.2, the OCSMM may be adopted to learn from data points whose uncertainties are observed explicitly. To illustrate this claim, we generate samples from the unit circle using $x = \cos \theta + \varepsilon$ and $y = \sin \theta + \varepsilon$ where $\theta \sim (-\pi, \pi]$ and ε is a zero-mean isotropic Gaussian noise $\mathcal{N}(0, 0.05)$. A different point-wise Gaussian noise $\mathcal{N}(0, \omega_i)$ where $\omega_i \in (0.2, 0.3)$ is further added to each point to simulate the random measurement corruption. In this experiment, we assume that ω_i is available during training. This situation is often encountered in many applications such as astronomy and computational biology. Both OCSVM and OCSMM are trained on the corrupted data. As opposed to the OCSVM that considers only the observed data points, the OCSMM also uses ω_i for every point via the kernel (7). Then, we consider a slightly more complicated data generated by $x = r \cdot \cos(\theta)$ and $y = r \cdot \sin(\theta)$ where $r = \sin(4\theta) + 2$ and $\theta \in (0, 2\pi]$. The data used in both examples are illustrated in Figure 4.

As illustrated by Figure 4, the density function estimated by the OCSMM is relatively less susceptible to the additional corruption than that estimated by the OCSVM, and tends to estimate the true density more accurately. This is not surprising because we also take into account an additional information about the uncertainty. However, this experiment suggests that when dealing with uncertain data, it might be ben-

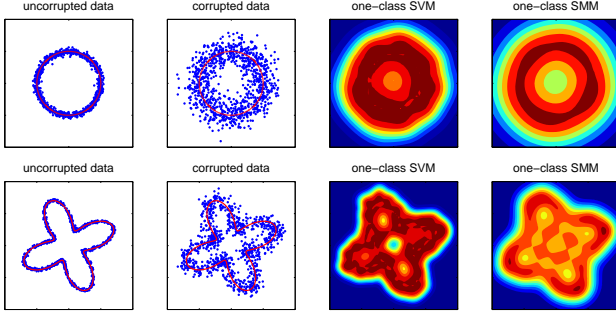


Figure 4: The density functions estimated by the OCSVM and the OCSMM using the corrupted data.

eficial to also estimate the uncertainty, as commonly performed in astronomy, and incorporate it into the model. This scenario has not been fully investigated in AI and machine learning communities. Our framework provides one possible way to deal with such a scenario.

5.3 Sloan Digital Sky Survey

Sloan Digital Sky Survey (SDSS)² consists of a series of massive spectroscopic surveys of the distant universe, the milky way galaxies, and extrasolar planetary systems. The SDSS datasets contain images and spectra of more than 930,000 galaxies and more than 120,000 quasars.

In this experiment, we are interested in identifying anomalous groups of galaxies, as previously studied in Póczos et al. (2011) and Xiong et al. (2011a;b). To replicate the experiments conducted in Xiong et al. (2011b), we use the same dataset which consists of 505 spatial clusters of galaxies. Each of which contains about 10-15 galaxies. The data were preprocessed by PCA to reduce the 1000-dimensional features to 4-dimensional vectors.

To evaluate the performance of different algorithms to detect group anomaly, we consider artificially random injections. Each anomalous group is constructed by randomly selecting galaxies. There are 50 anomalous groups of galaxies in total. Note that although these groups of galaxies contain usual galaxies, their aggregations are anomalous due to the way the groups are constructed.

The average precision (AP) and area under the ROC curve (AUC) from 10 random repetitions are shown in Figure 5. Based on the average precision, KNN-L2, MGM, and OCSMM achieve similar results on this dataset and KNN-Renyi outperforms all other algorithms. On the other hand, the OCSMM and KNN-

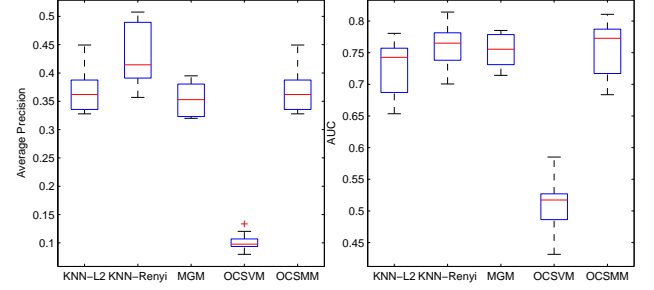


Figure 5: The average precision (AP) and area under the ROC curve (AUC) of different group anomaly detection algorithms on the SDSS dataset.

Renyi achieve highest AUC scores on this dataset. Moreover, it is clear that point anomaly detection using the OCSVM fails to detect group anomalies.

5.4 High Energy Particle Physics

In this section, we demonstrate our group anomaly detection algorithm in high energy particle physics, which is largely the study of fundamental particles, e.g., neutrinos, and their interactions. Essentially, all particles and their dynamics can be described by a quantum field theory called the *Standard Model*. Hence, given massive datasets from high-energy physics experiments, one is interested in discovering deviations from known Standard Model physics.

Searching for the Higgs boson, for example, has recently received much attention in particle physics and machine learning communities (see e.g., Bhat (2011), Vatanen et al. (2012) and references therein). A new physical phenomena usually manifest themselves as tiny excesses of certain types of collision events among a vast background of known physics in particle detectors.

Anomalies occur as a cluster among the background data. The background data distribution contaminated by these anomalies will therefore be different from the true background distribution. It is very difficult to detect this difference in general because the contamination can be considerably small. In this experiment, we consider similar condition as in Vatanen et al. (2012) and generate data using the standard HEP Monte Carlo generators such as PYTHIA³. In particular, we consider a Monte Carlo simulated events where the Higgs is produced in association with the W boson and decays into two bottom quarks.

The data vector consists of 5 variables (p_x, p_y, p_z, e, m) corresponding to different characteristics of the topology of a collision event. The variables p_x, p_y, p_z, e rep-

²See <http://www.sdss.org> for the detail of the surveys.

³<http://home.thep.lu.se/~torbjorn/Pythia.html>

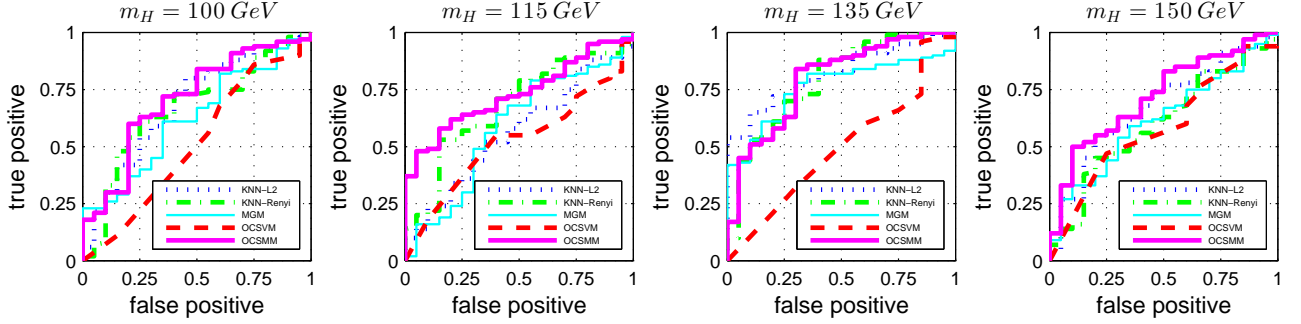


Figure 6: The ROC of different group anomaly detection algorithms on the Higgs boson datasets with various Higgs masses m_H . The associated AUC scores for different settings, sorted in the same order appeared in the figure, are (0.6835,0.6655,0.6350,0.5125,**0.7085**), (0.5645,0.6783,0.5860,0.5263,**0.7305**), (**0.8190**,0.7925,0.7630,0.4958,0.7950), and (0.6713,0.6027,0.6165,0.5862,**0.7200**).

resents the momentum four-vector in units of GeV with $c = 1$. The variable m is the particle mass in the same unit. The signal looks slightly different for different Higgs masses m_H , which is an unknown free parameter in the Standard Model. In this experiment, we consider $m_H = 100, 115, 135$, and 150 GeV. We generate 120 groups of collision events, 100 of which contain only background signals, whereas the rest also contain the Higgs boson collision events. For each group, the number of observable particles ranges from 200 to 500 particles. The goal is to detect the anomalous groups of signals which might contain the Higgs boson without prior knowledge of m_H .

Figure 6 depicts the ROC of different group anomaly detection algorithms. The OCSMM and KNN-based group anomaly detection algorithms tend to achieve competitive performance and outperform the MGM algorithm. Moreover, it is clear that traditional point anomaly detection algorithm fails to detect high-level anomalous structures.

6 Conclusions and Discussions

To conclude, we propose a simple and efficient algorithm for detecting group anomalies called one-class support measure machine (OCSMM). To handle aggregate behaviors of data points, groups are represented as probability distributions which account for higher-order information arising from those behaviors. The set of distributions are represented as mean functions in the RKHS via the kernel mean embedding. We also extend the relationship between the OCSVM and the KDE to the OCSMM in the context of variable kernel density estimation, bridging the gap between large-margin approach and kernel density estimation. We demonstrate the proposed algorithm on both synthetic and real-world datasets, which achieve competitive results compared to existing group anomaly de-

tection techniques.

It is vital to note the differences between the OCSMM and hierarchical probabilistic models such as MGM and FGM. Firstly, the probabilistic models assume that data are generated according to some parametric distributions, i.e., mixture of Gaussian, whereas the OCSMM is nonparametric in the sense that no assumption is made about the distributions. It is therefore applicable to a wider range of applications. Secondly, the probabilistic models follow a bottom-up approach. That is, detecting group-based anomalies requires point-based anomaly detection. Thus, the performance also depends on how well anomalous points can be detected. Furthermore, it is computationally expensive and may not be suitable for large-scale datasets. On the other hand, the OCSMM adopts the top-down approach by detecting the group-based anomalies directly. If one is interested in finding anomalous points, this can be done subsequently in a group-wise manner. As a result, the top-down approach is generally less computational expensive and can be used efficiently for online applications and large-scale datasets.

References

- I. S. Abramson. On bandwidth variation in kernel estimates-a square root law. *The Annals of Statistics*, 10(4):1217–1223, 1982.
- A. Berline and T. C. Agnan. *Reproducing Kernel Hilbert Spaces in Probability and Statistics*. Kluwer Academic Publishers, 2004.
- P. C. Bhat. Multivariate Analysis Methods in Particle Physics. *Ann.Rev.Nucl.Part.Sci.*, 61:281–309, 2011.
- D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. *Journal of Machine Learning Research*, 3: 993–1022, 2003.
- J. Bovy, J. F. Hennawi, D. W. Hogg, A. D. Myers, J. A. Kirkpatrick, D. J. Schlegel, N. P. Ross, E. S. Sheldon,

- I. D. McGreer, D. P. Schneider, and B. A. Weaver. Think outside the color box: Probabilistic target selection and the sdss-xdqs0 quasar targeting catalog. *The Astrophysical Journal*, 729(2):141, 2011.
- L. Breiman, W. Meisel, and E. Purcell. Variable kernel estimates of multivariate densities. *Technometrics*, 19(2):135–144, 1977.
- P. K. Chan and M. V. Mahoney. Modeling multiple time series for anomaly detection. In *Proceedings of the 5th IEEE International Conference on Data Mining (ICDM)*, pages 90–97. IEEE Computer Society, 2005.
- V. Chandola, A. Banerjee, and V. Kumar. Anomaly detection: A survey. *ACM Comput. Surv.*, 41(3):1–58, July 2009.
- K. Das, J. Schneider, and D. B. Neill. Anomaly pattern detection in categorical datasets. In *ACM-SIGKDD*, pages 169–176. ACM, 2008.
- K. Fukumizu, F. R. Bach, and M. I. Jordan. Dimensionality reduction for supervised learning with reproducing kernel hilbert spaces. *Journal of Machine Learning Research*, 5:73–99, December 2004.
- H. Hoffmann. Kernel PCA for novelty detection. *Pattern Recognition*, 40(3):863–874, 2007.
- J. A. Kirkpatrick, D. J. Schlegel, N. P. Ross, A. D. Myers, J. F. Hennawi, E. S. Sheldon, D. P. Schneider, and B. A. Weaver. A simple likelihood method for quasar target selection. *The Astrophysical Journal*, 743(2):125, 2011.
- K. Muandet, K. Fukumizu, F. Dinuzzo, and B. Schölkopf. Learning from distributions via support measure machines. In *Advances in Neural Information Processing Systems (NIPS)*, pages 10–18. 2012.
- B. Póczos, L. Xiong, and J. G. Schneider. Nonparametric divergence estimation with applications to machine learning on distributions. In *Proceedings of the 27th Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 599–608, 2011.
- N. P. Ross, A. D. Myers, E. S. Sheldon, C. Yche, M. A. Strauss, J. Bovy, J. A. Kirkpatrick, G. T. Richards, R. Aubourg, M. R. Blanton, W. N. Brandt, W. C. Carithers, R. A. C. Croft, R. da Silva, K. Dawson, D. J. Eisenstein, J. F. Hennawi, S. Ho, D. W. Hogg, K.-G. Lee, B. Lundgren, R. G. McMahon, J. Miralda-Escud, N. Palanque-Delabrouille, I. Pris, P. Petitjean, M. M. Pieri, J. Rich, N. A. Roe, D. Schiminovich, D. J. Schlegel, D. P. Schneider, A. Slosar, N. Suzuki, J. L. Tinker, D. H. Weinberg, A. Weyant, M. White, and W. M. Wood-Vasey. The sdss-iii baryon oscillation spectroscopic survey: Quasar target selection for data release nine. *The Astrophysical Journal Supplement Series*, 199(1):3, 2012.
- B. Schölkopf and A. J. Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, Cambridge, MA, USA, 2001. ISBN 0262194759.
- B. Schölkopf, J. C. Platt, J. C. Shawe-Taylor, A. J. Smola, and R. C. Williamson. Estimating the support of a high-dimensional distribution. *Neural Computation*, 13:1443–1471, 2001.
- A. Smola, A. Gretton, L. Song, and B. Schölkopf. A hilbert space embedding for distributions. In *Proceedings of the 18th International Conference on Algorithmic Learning Theory*, pages 13–31. Springer-Verlag, 2007.
- B. K. Sriperumbudur, A. Gretton, K. Fukumizu, B. Schölkopf, and G. R. G. Lanckriet. Hilbert space embeddings and metrics on probability measures. *Journal of Machine Learning Research*, 2010.
- D. M. J. Tax and R. P. W. Duin. Support vector domain description. *Pattern Recognition Letters*, 20:1191–1199, 1999.
- D. M. J. Tax and R. P. W. Duin. Support vector data description. *Machine Learning*, 54(1):45–66, 2004.
- G. R. Terrell and D. W. Scott. Variable kernel density estimation. *The Annals of Statistics*, 20(3):1236–1265, 1992.
- T. Vatanen, M. Kuusela, E. Malmi, T. Raiko, T. Aaltonen, and Y. Nagai. Semi-supervised detection of collective anomalies with an application in high energy particle physics. In *IJCNN*, pages 1–8. IEEE, 2012.
- L. Xiong, B. Póczos, and J. Schneider. Group anomaly detection using flexible genre models. In *Proceedings of Advances in Neural Information Processing Systems (NIPS)*, 2011a.
- L. Xiong, B. Póczos, J. G. Schneider, A. Connolly, and J. VanderPlas. Hierarchical probabilistic models for group anomaly detection. *Journal of Machine Learning Research - Proceedings Track*, 15:789–797, 2011b.
- M. Zhao and V. Saligrama. Anomaly detection with score functions based on nearest neighbor graphs. In *Proceedings of Advances in Neural Information Processing Systems (NIPS)*, pages 2250–2258, 2009.